

Application-centric Resource Provisioning for Amazon EC2 Spot Instances

Sunirmal Khatua ^{#1}, Nandini Mukherjee ^{§2}

[#] *Department of Computer Science and Engineering, University of Calcutta, India*

¹ skhatuacomp@caluniv.ac.in

[§] *Department of Computer Science and Engineering, Jadavpur University, India*

² nmukherjee@cse.jdvu.ac.in

Abstract—In late 2009, Amazon introduced spot instances to offer their unused resources at lower cost with reduced reliability. Amazon's spot instances allow customers to bid on unused Amazon EC2 capacity and run those instances for as long as their bid exceeds the current spot price. The spot price changes periodically based on supply and demand, and customers whose bids exceed it gain access to the available spot instances. Customers may expect their services at lower cost with spot instances compared to on-demand or reserved. However the reliability is compromised since the instances(IaaS) providing the service(SaaS) may become unavailable at any time without any notice to the customer. Checkpointing and migration schemes are of great use to cope with such situation. In this paper we study various checkpointing schemes that can be used with spot instances. Also we devise some algorithms for checkpointing scheme on top of application-centric resource provisioning framework that increase the reliability while reducing the cost significantly.

I. INTRODUCTION

The era of cloud computing provides high utilization and high flexibility of managing the computing resources. The elasticity and on demand availability features of cloud computing ensure high utilization of resources. Furthermore, resources can be availed from templates that enforce standards so that resources can be used with best management considerations without prior knowledge. Therefore, flexibility is also high in cloud environment. The cloud computing service models incorporate Infrastructure as a Service (IaaS), Platform as a Service (PaaS) and Software as a Service (SaaS). IaaS provides raw computing resources with different capacity in the form of Virtual Machines (VM). Cloud service providers, like Google [13], Amazon [12] provide these services and charge prices against these services from the clients. Among many such providers, Amazon defines the capacity of resources in the form of 64 instance types [9] based on storage, compute unit and I/O performance. The cost of these instance types depends on the purchasing models defined by Amazon namely on-demand, reserved and spot. On-Demand Instances let one pay for compute capacity by the hour with no long-term commitments or upfront payments. However with On-Demand Instances one may not have access to the resources immediately. On the other hand, Reserved Instances facilitate the client to make a low, one-time, upfront payment for an instance, reserve it and get significant discount on hourly charge over On-Demand Instances. Reserved Instances are

always available for the durations for which the clients reserve. In contrast with the above two policies, where rates are fixed, Spot Instances provide the ability for customers to purchase compute capacity with no upfront commitment and at a variable hourly rates with a customer-defined upper bound (bid) on the rate. Spot Instances are available only during the time when the spot price is below the customer defined bid.

Thus spot instances make the resources unreliable in nature and inappropriate for long running jobs like image processing, gene sequence analysis etc. At the same time they offer the opportunity to accomplish such jobs at a much lower cost than on demand or reserved policies. Clearly checkpointing may be a good option to make a tradeoff between the cost and reliability. Checkpointing allows to store a snapshot of the current application state, and later on, use it for restarting the execution at an opportunistic moment.

Various checkpointing techniques have been discussed in [3] to provide reliability with Amazon spot instances at lower cost. In this paper we study some of these techniques and evaluate their performances. We also investigate the effectiveness of application centric resource provisioning framework [2] for actively monitoring the deployed spot instances for an application and for taking necessary actions as the spot instances become unavailable or the spot price changes. Finally we propose and evaluate a novel checkpointing scheme for the application centric resource provisioning framework.

The rest of the paper is organized as follows. A brief review of the related works is presented in Section II. An overview of the application centric resource provisioning framework is given in Section III. The available resource provisioning options are described in Section IV. Section V discusses the existing checkpointing schemes for spot instances while a proposed checkpointing scheme for the application centric resource provisioning framework is described in Section VI. A simulated result for comparing the proposed checkpointing scheme with existing ones is presented in Section VII. Finally we conclude with a direction of future work in Sections VIII.

II. RELATED WORK

During the last couple of years, a lot of works [2][5]-[7] concentrate on the cloud management aspect from the economic point of view. Most of them adapt a middleware

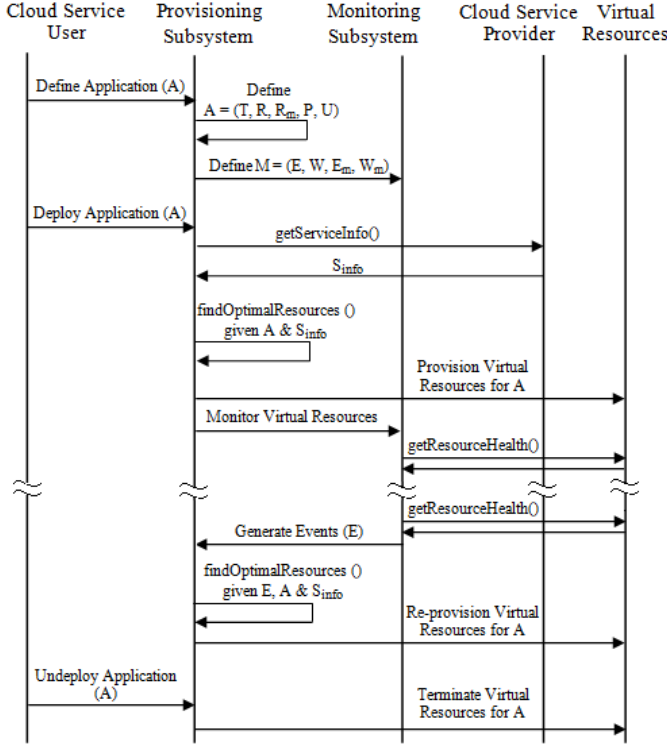


Fig. 2. Resource provisioning algorithm

$$M = (E, W, E_m, W_m) \quad (2)$$

where E is the set of events, $\{e\}$
 W is the set of workflows $\{w\}$
 $E_m : E \rightarrow T \mid E \rightarrow R$
 $W_m : W \rightarrow E$

A brief description of functioning of the application centric resource provisioning framework is depicted in Figures 1 and 2. The CSUs use a Unified Client API to define her application according to equations 1 and 2. This unified definition of the applications are used by two important subsystems of the framework namely provisioning subsystem and monitoring subsystem as described below.

A. Provisioning Subsystem

The provisioning subsystem determines optimal provisioning of virtual resources for an application(A) satisfying the policies(P) specified for it. The application's required service level is stored in the policy(P). The provisioning subsystem queries various providers to get information about their offered services(S_{info}). S_{info} consists of provider id, service id, QoS id and the associated cost. The provisioning subsystem uses P(desired service level), S_{info} and an optimization algorithm to find the optimal resource requirement for the application while maintaining the desired service level.

B. Monitoring Subsystem

The monitoring subsystem implements a feedback system to inform the provisioning subsystem about the current state of

the deployed application. The monitoring subsystem actively monitors the state of the deployed application and generates various events [1] to designate a change in the state. In the proposed framework, an application can be in any of the six defined states, namely New, Inactive, Active, Unbalanced, Unreachable and Terminated (Figure 3). Initially any application is in the New state. Once such an application is mapped to various modules according to the unified definition, the application enters into the Inactive state. In the Inactive state the application is composed (as specified by the unified definition), the required infrastructure is programmed (as determined by the optimization algorithm) and is ready to be deployed within the cloud. The application is then deployed and becomes ready to be accessed via the corresponding URL and its state changes to the Active state. When in the Active state, the user pays for the cloud resources. An application can be moved to the inactive state or to the active state manually for fine tuning. If the application is no longer required, the user can release the mapping from the middleware and the application will be in the Terminated state.

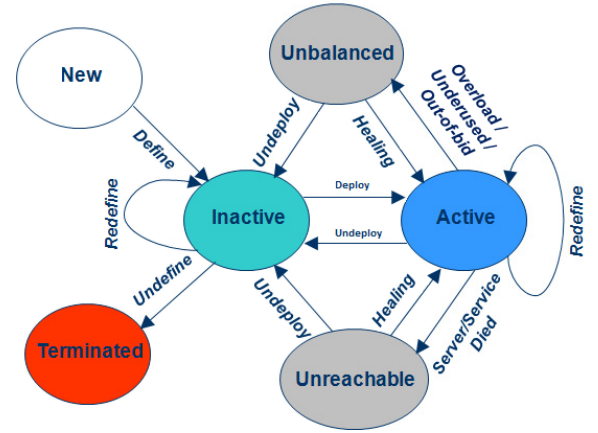


Fig. 3. States of an Application in application-centric cloud

Two other important states of the application are the Unbalanced and Unreachable states. If the deployed application is overloaded or underused based on certain threshold conditions, it reaches the Unbalanced state. Similarly if any of the resource which is deployed for the application fails, the application goes into the Unreachable state. In these states, a workflow can be maintained or generated in order to heal the situation and these actions send back the application to the Active state. Figure 3 depicts the states of the application.

The monitoring subsystem uses different event generation schemes for its proper operation. In [2], five event generation schemes are described. These are threshold based, prediction based, request based, ping based and schedule based. The generated events carry necessary information needed to re-provision the application resources to optimize or heal some undesired situation. Once an event is generated, the monitoring subsystem sends the event to the provisioning subsystem. Once an event(E) is received, the provisioning subsystem analyzes the event and uses E, P, S_{info} and an optimization algorithm

for reprovisioning the application onto appropriate resources.

IV. RESOURCE PROVISIONING IN AMAZON EC2 CLOUD

In this paper multiple providers of application centric resource provisioning are not considered. Rather, we consider various resource provisioning options available from Amazon EC2 public provider only. Amazon sells their resources in the form of on-demand, reserved and spot instances.

On-demand resources can be used without any upfront payment and just paying as much as the client use on a hourly basis. However request for on-demand instances may not be met immediately due to unavailability of Amazon EC2 resources. Thus for a long term and time critical application it is required to opt for reserved instances. With reserved instances required resources can be reserved with some upfront payment and access to the reserved instances can be made whenever the client needs. Amazon also provides competitive discounts on the hourly charge for the reserved instances. The third category of the instances, i.e. spot instances, allow the user to use Amazon's unused resources at lower cost compared to on-demand and reserved instances if available. The prices of spot instances, called spot price, depends on the demand and supply of the specific instance type at a specific availability zone. Users need to define the bid (the maximum cost he is willing to pay per instance) for a specific instance type at a specific availability zone and the spot instance request will be granted if the current spot price is less than the bid defined by the user.

Characteristics of Amazon EC2 Spot Instances:

The variable price of spot instances makes them an important consideration for optimizing resource requirement for an application. However, their volatile nature makes them inherently unreliable and hence the optimization algorithms become more challenging than the other instances.

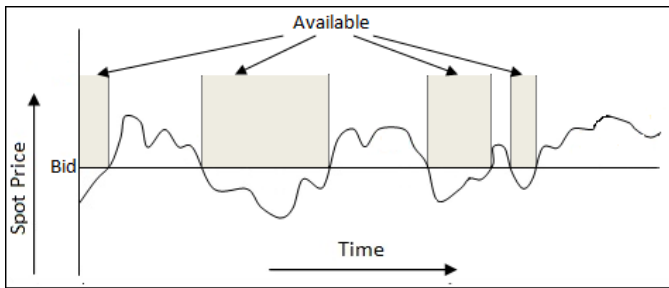


Fig. 4. Availability of a spot instance

Various characteristics of Amazon EC2 spot instances [10] are summarized below:

- Spot instances are available when the user's bid exceeds the current spot price (Fig. 4).
- Spot instances are terminated (becomes unavailable) without any notification to the user whenever the current spot price exceeds the user's bid.
- The price per instance-hour for a spot instance is set at the beginning of each instance-hour. Any change to the

spot price will not be reflected until the next instance-hour begins.

- Amazon will not charge the last partial hour if the spot instance is terminated due to out-of-bid situation. However Amazon will charge the full hour if the user terminate the instance forcefully.
- Amazon provides the history of spot prices of a spot instance at a specific availability zone for the last 3 months free of cost.

V. CHECKPOINTING SCHEMES FOR SPOT INSTANCES

The characteristics of spot instances make them appealing for long running jobs with divisible workloads [8]. Various existing checkpointing schemes can be adopted for saving the completed tasks and resuming the remaining tasks as and when the spot instances become available.

Existing Checkpointing Schemes:

The checkpointing schemes proposed in [3] are briefly described below:

1. No Checkpointing (NONE): Checkpoints are not taken and all the tasks for a job are required to be repeated after every out-of-bid events.
2. Optimal Checkpointing (OPT): Checkpoints are taken just prior to the out-of-bid events. Clearly it will save the maximum number of tasks out of each available interval for a given instance type and a user's bid.
3. Hourly Checkpointing (HOUR): Checkpoints are taken just prior to the beginning of next instance hour. Since Amazon is not charging any partial hour, this scheme will save as much tasks as the user is paying.
4. Rising edge-driven Checkpointing (EDGE): Checkpoints are taken after every increase (rising edge) of the current spot price.
5. Adaptive Checkpointing (ADAPT): Checkpoints are taken or skipped at regular intervals based on the expected recovery time for taking or skipping the checkpoint. It will take a checkpoint if the expected recovery time is higher for skipping the checkpoint. The expected recovery time is calculated using a probability density function of expected out-of-bid events. Such a probability density function is determined from the history of spot prices and the user defined bid.

Out of the above five checkpointing schemes NONE and OPT provide two extreme results without any practical value. They are used to provide comparative study of the other realistic checkpointing schemes.

VI. A NOVEL CHECKPOINTING SCHEME FOR APPLICATION-CENTRIC RESOURCE PROVISIONING

In this section we propose a novel checkpointing scheme for spot instances on top of application-centric resource provisioning framework. For the purpose we devise a new event generation scheme that deals with spot instances. The new checkpointing scheme is targeted to achieve performance comparative to OPT checkpointing scheme described above. Before describing the scheme, we introduce a modified event

generation scheme for our application-centric resource provisioning framework.

A. Event Generation Scheme for Spot Instances

The event generation schemes proposed in [2] is extended to include new events that support spot instances. As discussed in Section IV, the availability of spot instances depends on the current spot price and the user defined bid. Also spot instances become unavailable without prior notification to the clients that makes them inherently unreliable. The reliability can be increased by taking checkpoints (saving completed tasks) during the available periods. However, the time of taking checkpoints affects the reliability as well as job completion time and cost.

Accordingly, in this paper we propose a new event generation scheme to handle spot instances. Three events are proposed, namely E_{ckpt} , $E_{terminate}$ and E_{launch} . E_{ckpt} is used for taking checkpoint, $E_{terminate}$ is used to terminate a spot instance forcefully and E_{launch} is used to relaunch a previously terminated spot instance. We define two bid values for the purpose - one for the application (A_{bid}) and other for the spot instance (S_{bid}). S_{bid} is sufficiently large and is used in the request for spot instance. Clearly, the value is maintained at such a high level, that Amazon will never terminate the spot instances due to out-of-bid situation. On the other hand, A_{bid} is user defined bid for the application and is stored in the monitoring subsystem as part of the event definition (E) of the application. The Monitor module actively monitors the current spot price and generates the two events, E_{ckpt} and $E_{terminate}$, for the Controller module. On the basis of these two events, the Controller module either takes a checkpoint or terminate the corresponding spot instance respectively. However to increase the performance, the Controller module will query the current spot price only at specific points of time called decision points. Since the cost of spot instance is not changed during an instance hour and is fixed at the beginning of that instance hour, the decision points should be relative to the beginning of an instance hour. Accordingly we define two decision points just prior to each hour boundary as follows:

$$t_{cd} = t_h - t_c - t_w \quad (3)$$

$$t_{td} = t_h - t_w \quad (4)$$

where t_{cd} and t_{td} are the decision points for checkpointing and terminating a spot instance. t_h is an hour boundary, t_c is the time needed to take a checkpoint and t_w is the waiting time to get the current spot price. The Monitor module will generate E_{ckpt} at t_{cd} if the current spot price exceeds A_{bid} and will generate $E_{terminate}$ at t_{td} if the current spot price is still above the A_{bid} . It will generate E_{launch} at the start of each available period of a spot instance with respect to A_{bid} . This event generation scheme is illustrated in Figure 5. It will generate neither E_{ckpt} nor $E_{terminate}$ for the hour boundary t_{h1} . It will generate E_{ckpt} but not $E_{terminate}$ for the hour boundary t_{h2} . For the hour boundary t_{h3} , it will generate both E_{ckpt} and $E_{terminate}$ since the user will have to pay above A_{bid} for the next hour.

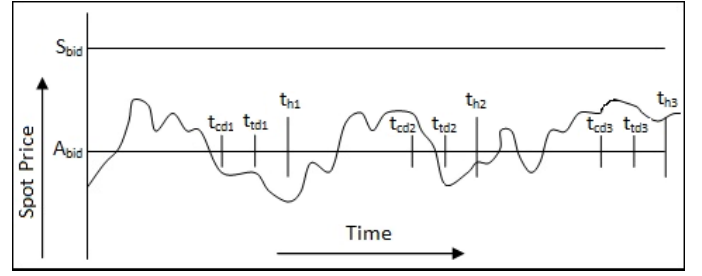


Fig. 5. Decision Points for Event Generation

B. The Application-Centric Checkpointing Scheme

In this section, we propose a checkpointing scheme on top of the application centric resource provisioning framework, called Application Centric Checkpointing (ACC). ACC is based on the event generation scheme discussed in the previous subsection and is described by the sequence diagram shown in Figure 6.

The following unified definition can be used for an application with divisible workloads to be run on spot:

$$A = (T, R, R_m, P, U, M) \quad (5)$$

where $T = \{t_1\}$

$R = \{r_1, r_2\}$

$r_1.provider = ec2, r_1.type = \text{spot instance},$
 $r_1.size = < instance_type >$

$r_2.provider = ec2, r_2.type = \text{EBS},$
 $r_2.size = 1\text{GB}$

$R_m = \{r_1 \rightarrow t_1, r_2 \rightarrow t_1\}$

$P = \{sla\}$

$$M = (E, W, E_m, W_m) \quad (6)$$

where $E = \{E_{ckpt}, E_{terminate}, E_{launch}\},$

$threshold \text{ for all events } = < A_{bid} > ,$

$E_{launch}.bid = < S_{bid} >$

$W = \{W_{start}, W_{ckpt}, W_{terminate}, W_{launch}\}$

$W_{start} = \{ \text{Launch spot};$
 $\text{Mount EBS};$
 $\text{Copy job to EBS};$
 $\text{Start job} \},$

$W_{ckpt} = \{ \text{Save results to EBS},$
 $W_{terminate} = \{ \text{Terminate spot} \} \ \&$

$W_{launch} = \{ \text{Launch spot};$
 $\text{Mount EBS};$
 $\text{Resume tasks} \},$

$E_m = \{E_{ckpt} \rightarrow r_1, E_{terminate} \rightarrow r_1,$
 $E_{launch} \rightarrow r_1\}$

$W_m = \{W_{ckpt} \rightarrow E_{ckpt},$
 $W_{terminate} \rightarrow E_{terminate},$
 $W_{launch} \rightarrow E_{launch}\}$

The Elastic Block Storage (EBS) [11] is used to save the completed tasks during checkpoint. The parameters

$instance_type$, A_{bid} and S_{bid} can be set either manually by the end user or by some optimization or greedy algorithms. The provisioning subsystem (Deployer module) can use the following simple greedy strategy for choosing A_{bid} and $instance_type$:

Algorithm 1 Determine A_{bid} & $instance_type$

1. Retrieve S_{info} from Amazon EC2.
/* S_{info} carries availability zone, spot instance type and history of spot price.*/
2. Find the list of instance types that meet the required service level agreement(sla) specified in P.
/* The list is denoted by L.*/
3. Calculate application bid as

$$A_{bid} = \min C_i, \forall i \in L \quad (7)$$

/* C_i is the corresponding on demand instance's cost per hour for the instance type i.*/

4. For each instance type $i \in L$
 - 4.1 Calculate Expected Execution Time (EET) for a job of length 'w' when executed in a spot instance of instance type 'i' with a bid of value A_{bid} .

$$EET_i = \frac{w \sum_{k=w}^{\infty} f_i(k) + \sum_{k=0}^{w-1} (k+r) f_i(k)}{1 - \sum_{k=0}^{w-1} f_i(k)} \quad (8)$$

/* $f_i(t)$ is the probability density function of the spot instance type i's failure for out-of-bid. The $f_i(t)$ is calculated from the spot instance type i's history of spot price and A_{bid} .*/

5. Choose $instance_type = i \mid EET_i$ is minimum.
-

After determining the parameters A_{bid} & $instance_type$, the Deployer module starts W_{start} workflow. The W_{start} workflow launches a spot instance as per the specification of the resource r_1 and an EBS volume as per the specification of the resource r_2 . The workflow then mounts the EBS volume to the spot instance, copy the job from the application repository to the EBS and starts the job.

Once the application is deployed, EC2 starts charging for the resources. The monitoring subsystem (Monitor module) calculates t_{cd} and t_{td} as per Equ. 3 & 4 for the current hour boundary. At t_{cd} the monitor module retrieves the current spot price(P). If P exceeds A_{bid} , it generates E_{ckpt} event for the Controller module. On receiving E_{ckpt} event, the Controller module executes W_{ckpt} workflow. The W_{ckpt} workflow just saves the results(the completed tasks) to the EBS volume. The Monitor module also retrieves the current spot price(P) at t_{td} . If P still exceeds A_{bid} , it generates $E_{terminate}$ event for the Controller module. On receiving $E_{terminate}$ event, the Controller module executes $W_{terminate}$ workflow. The $W_{terminate}$ workflow terminates the spot instance forcefully.

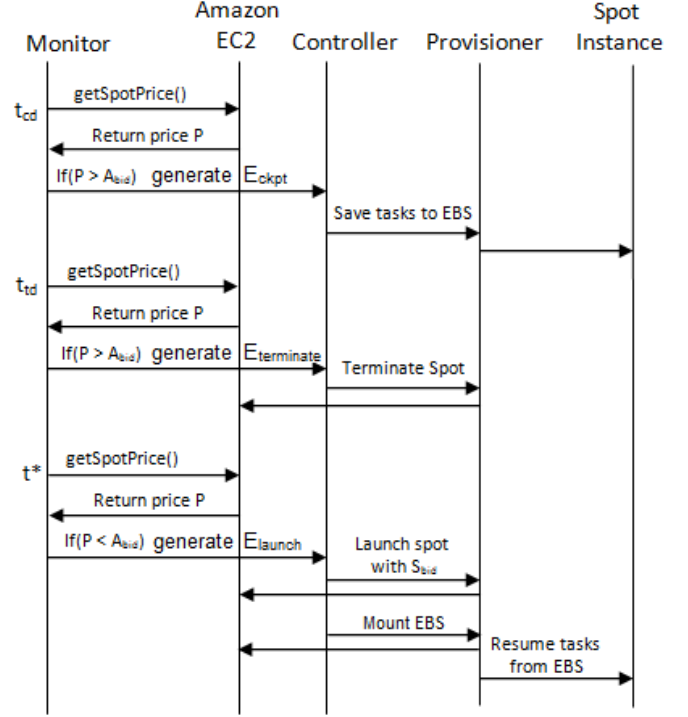


Fig. 6. Application Centric Checkpointing Scheme

The Monitor module repeats the above procedure till P does not exceed A_{bid} at t_{td} for all the subsequent hour boundaries.

If the instance is terminated at some t_{td} , the Monitor module will have to query for the current spot price to determine the next Available period(refer to Fig. 4) at some specific instance of time(t^*). However, the frequency of making the query is defined by the end user which may affect the job completion time slightly. At the start of the new available duration, the Monitor module generates E_{launch} event for the Controller module. On receiving E_{launch} event, the Controller module executes W_{launch} workflow. The W_{launch} workflow launches a new spot instance as specified in r_1 , mount the existing EBS volume to that instance and resume the remaining tasks of the job.

VII. IMPLEMENTATION AND EVALUATION

In this section we analyze and compare our proposed ACC checkpointing scheme with the existing checkpointing schemes. The experiments have been carried out on 64 spot instance types of Amazon EC2 those have also been used in [3]. The metrics used for this purpose include *job completion time*, *total monetary cost* and the *product of monetary cost x completion time* as the basis for comparison.

A. Simulation Setup

We have simulated the checkpointing schemes, discussed in section V & VI, using the same data set, parameters, algorithms and assumptions used in [3]. We have downloaded

the simulator [23] and applied the following modifications for our simulation setup:

- Modification is applied to all the checkpointing functions to rectify their [3] wrong assumption that Amazon charges each hour by the last price. The modified algorithm charges a spot instance by the cost of its instance type at the beginning of an instance-hour as specified in the characteristics of Amazon EC2 spot instances.
- A function is added to simulate the ACC checkpointing scheme discussed in Section VI-B.

In this paper we have not simulated the algorithm for determining A_{bid} and $instance_type$. Instead we have simulated the checkpointing schemes on all the 64 instance types under different A_{bid} values from \$0 to \$2 with a granularity of \$0.001.

B. Results and Discussion

We obtain the simulation result for *job completion time*, *total monetary cost* and the *product of monetary cost x completion time* for all the EC2 instance types. To simplify the discussion, we present the result of a linux based extra large (m1.xlarge) instance type in the eu-west-1 region. We concentrate on the performance of our proposed ACC checkpointing scheme compared to the optimal checkpointing scheme, OPT. We also include NONE, HOUR, EDGE and ADAPT checkpointing schemes in our result for completeness.



Fig. 7. Total monetary cost of Job completion

Fig 7 shows the comparison of total monetary cost needed to complete a job of length 500 minutes under different user's bid (A_{bid}) from \$0.401 to \$0.441. The result shows that ACC reduces the job completion cost significantly over the other realistic checkpointing schemes. However the cost is increased by 5.94% on average (min 0.33%, max 10.30%) compared to OPT scheme. This is because the OPT scheme can execute some fraction of the job free of cost for the partial hours.

In Fig. 8 we illustrate the comparison of various checkpointing schemes for the metric *job completion time*. Here we observe that ACC scheme outperforms all the checkpointing schemes including OPT. This is because ACC allows the job to continue even when the current spot price exceeds A_{bid} (in between a t_{td} and the corresponding hour boundary) without

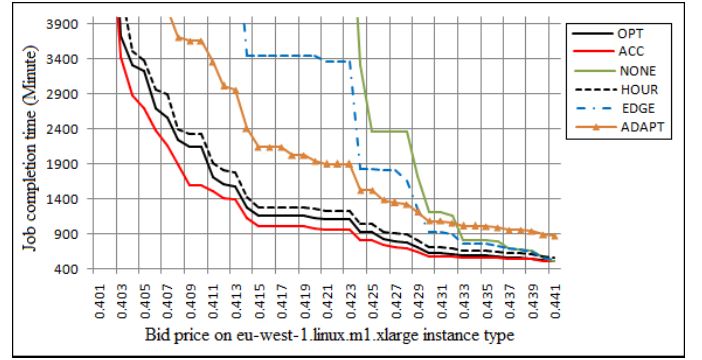


Fig. 8. Job completion time

affecting the job completion cost. The ACC scheme reduces the job completion time by an average value of 10.77% over the OPT scheme.

We plot the comparative study for the *product of monetary cost x completion time* in Fig. 9. Here also we observe that the ACC scheme reduce this metric by an average value of 5.56% over the OPT scheme.

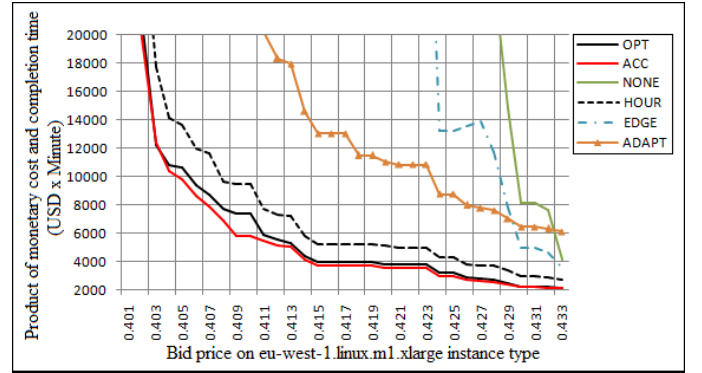


Fig. 9. Product of total cost and completion time

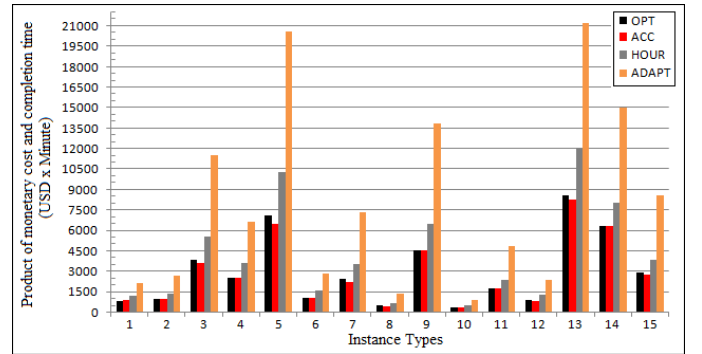


Fig. 10. Product of cost and completion time for different instance types

To gain confidence in our result, we have computed the average values of the above mentioned metrics for different bid values on all the 64 instance types. A sample of 15 difference instance types for the metric *product of monetary cost x completion time* is shown in

Fig. 10. For these 15 instance types, a gain of 4.03% for ACC over OPT is observed. We also observe that such percentage gain is increased for costly instance types.

In the previous research work [3], the authors conclude that OPT is the optimal checkpointing scheme and none of the practical schemes can perform better than OPT. That is true only if we use the same bid values for launching the spot instance and computing the checkpoint. However our proposed ACC checkpointing scheme perform very close to OPT or even better than OPT by separating these two bid values. Thus ACC outperforms all the existing checkpointing schemes for spot instances till date.

VIII. CONCLUSION AND FUTURE WORK

Checkpointing plays an important role in reliability of job execution over EC2 spot instances. In this paper we propose a checkpointing scheme on top of application-centric resource provisioning framework that not only increase the reliability but also reduces the cost significantly over the existing checkpointing schemes. The job completion cost under the proposed scheme is very close to the optimal checkpointing scheme. Even it performs better than the optimal scheme from the point of view of job completion time, as well as product of job completion time and cost.

In future we want to investigate more on the following issues:

- What is the optimal bid and the corresponding instance type for a given job?
- Should we migrate to another instance type during un-available period?
- What should be the new bid and the corresponding instance type for the migration?

REFERENCES

- [1] S. Khatua, A. Ghosh and N. Mukherjee, *Optimizing the utilization of virtual resources in Cloud environment*. In Proceedings of the VECIMS 2010, pp 82-87, DOI:10.1109/VECIMS.2010.5609349
- [2] S. Khatua, A. Ghosh and N. Mukherjee, *Application-centric Cloud managemengt*. In 9th IEEE/ACS International Conference on Computer Systems and Applications (AICCSA) 2011, pp 9-15, DOI:10.1109/AICCSA.2011.6126627
- [3] S. Yi, A. Andrzejak, D. Kondo, *Monetary Cost-Aware Checkpointing and Migration on Amazon Cloud Spot Instances*. In IEEE Transactions on Services Computing 2011. Volume: PP. Issue: 99.
- [4] P. Padala et al, *Adaptive control of virtualized resources in utility computing enironments*. In Proceedings of EuroSys, 2007.
- [5] Q. Li and Y. Guo, *Optimization of Resource Scheduling in Cloud Computing*. In 12th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing 2010, pp 315-320
- [6] S. Pandey, L. Wu, S. Mayura Guru and R. Buyya, *A Particle Swarm Optimization-based Heuristic for Scheduling Workflow Applications in Cloud Computing Environments*. In 24th IEEE International Conference on Advanced Information Networking and Applications 2010, pp 400-407
- [7] S. Chaisiri, B. Lee and D. Niyato, *Optimization of Resource Provisioning Cost in Cloud Computing*. In IEEE Transactions on Services Computing 2011
- [8] Y. Yang and H. Casanova, *Umr: A multi-round algorithm for scheduling divisible workloads*. In IPDPS, 2003, p 24.
- [9] Amazon EC2 Instance Types. Available from: <http://aws.amazon.com/ec2/instance-types/>
- [10] Amazon EC2 spot instances. Available from: <http://aws.amazon.com/ec2/spot-instances/>
- [11] Elastic Block Storage. Available from: <http://aws.amazon.com/ec2/ebs/>
- [12] S. Garfinkel, *An Evaluation of Amazons Grid Computing Services: EC2, S3 and SQS*. Tech. Rep. TR-08-07, Harvard University, August 2007.
- [13] Google Cloud Offering. Available from: <http://cloud.google.com/products/>
- [14] P. Barham, et al, *Xen and the Art of Virtualization*. In Proceedings of the 19th ACM symposium on Operating Systems Principles, 2003.
- [15] R. Wolsky et al, *Eucalyptus: A Technical Report on an Elastic Utility Computing Archietcture Linking Your Programs to Useful Systems*. Tech. Rep. 2008-10, University of California, Santa Barbara, October 2008.
- [16] Zabbix : an enterprise-class open source distributed monitoring solution for networks and applications. Available from: <http://www.zabbix.com/>
- [17] T. Harmer, et al, *An application-centric model for cloud management*. In Proceedings of 6th World Congress on Services, July 2010, pp 439-446, DOI:10.1109/SERVICES.2010.132
- [18] H.C.Lim et al, *Automated control in cloud computing: challenges and opportunities*. In Proceedings of the 1st workshop on Automated control for datacenters and clouds, Spain, June 19, 2009.
- [19] Mills and C. Terence, *Time Series Techniques for Economists*. Cambridge University Press, 1990.
- [20] Rajkumar Buyya, et al, *Market-Oriented Cloud Computing: Vision, Hype, and Reality for Delivering IT Services as Computing Utilities*. In 10th IEEE International Conference on High Performance Computing and Communications, 2008 (HPCC 08), pp.5-13, September 2008.
- [21] W. Iqbal, M. N. Dailey, D. Carrera, and P. Janecek, *Adaptive resource provisioning for read intensive multi-tier applications in the cloud*. In Future Generation of Computer Systems, Vol. 27 Issue 6, pp.871-879, June 2011.
- [22] Jin Shao and Qianxiang Wang, *A Performance Guarantee Approach for Cloud Applications Based on Monitoring*. In 35th IEEE Annual Computer Software and Applications Conference Workshops (COMP-SACW), 2011, pp.25-30, July 2011
- [23] Checkpointing Simulator for spot instances. Available from: <http://spotckpt.sourceforge.net>